

# AI Math Agents

## LLMs, Information Systems Biology, Precision Health, and Alzheimer's Genomics

Melanie Swan,\* DIYgenomics, University College London (Research Associate CBT); Takashi Kido, Advanced Comprehensive Research Organization, Teikyo University, Preferred Networks, Inc.; Eric Roland, RedBud AI, LLC; Renato P. dos Santos, Centre for Generative AI in Cognition and Education, Lutheran University of Brazil

### Background

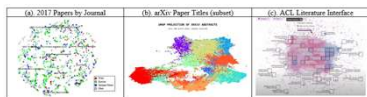
Generative AI technologies may be applied to advance science in new ways. LLMs (large language models) are computerized language models generated primarily with transformer neural networks (deep learning methods) having billions of parameters, and pre-trained on large data corpora, e.g. GPT-4 (OpenAI), LaMDA (Google), and LLaMA (Meta AI). Transformer neural networks are an advance which allows a whole data corpus to be processed simultaneously to assess connections between data elements. LLMs are seen as being crucial for genomic analysis given the large scope of data and interrelated parts [1,2].

The important result of LLMs is that they are a linguistic user interface, a language-based access tool, via natural language for human-AI chat (familiar from chatGPT), but more extensively, via formal languages such as programmatic code and mathematics, for further build-out of the computational infrastructure. The entirety of data corpora, not only word-based knowledge bases but also software code and mathematics, are being digitized and mobilized as available easy-to-use tools. Into this trajectory, the current work formulates an AI-based mathematics approach to Alzheimer's genomics, introducing Math Agents, mathematical embedding, and equation clusters as tools for representing and possibly evaluating mathematical ecologies/mathscapes (sets of equations).

### Method

AI-enabled mathematics tools are introduced and demonstrated in the case of Alzheimer's disease and aging as information systems biology problems. A theoretical model is elaborated, applying multiscalar physics mathematics (elucidating near-far entropic correlations in systems) to disease mathematics and whole-human genomic data for two precision medicine participants. Vector embedding as a standard machine learning method is employed with mathematical equations and genomic data as the input.

Figure 1. Embedding Visualization examples with Academic Papers as the Data Corpus



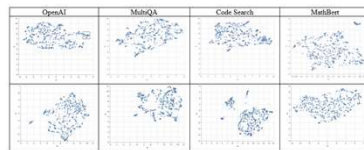
Vector embedding is the algorithmic processing of data into character strings for high dimensional analysis with results translated back into low dimensional (2D) output for interpretation. Used routinely in machine learning, the method primarily targets traditional content types, namely, text, images, and sound. Word embeddings are used in LLMs to process news, social media, and other online data corpora such as Wikipedia. Embeddings are starting to be developed for the analysis of scientific content, however, still targeted to words (Figure 1). Visualizations of embeddings for three projects, each of which uses the entirety of an academic literature as the data corpus, are presented: (a) all papers published in 2017 by journal [3], (b) all arXiv paper titles and abstracts (2.3 million) [4], and (c) an interactive visualization of the ACL Anthology (Association of Computational Linguistics) of 85,000 papers [5]. The interpretation is that clustering provides a relevant signal, with zoom-explorer functionality as visibility into the data set.

### Mathematical Embedding

The mathematical embedding is novel as an equation encoded in the form of a character-based vector string for input to high dimensional analysis in machine learning systems. The research aim is implementing mathematics as an approach to solving problems in data-intensive areas such as whole-human genome analysis and Alzheimer's disease. Causality is difficult to trace across tiers of complexity in biosystems. Hence, mathematics for modeling multiscalar physical systems is indicated, particularly renormalization mathematics which formalize a system-wide factor such as symmetry (in the universe) or free energy (in biological systems) conserved across scale tiers. Two recognized renormalization methods are AdS/CFT and Chern-Simons.

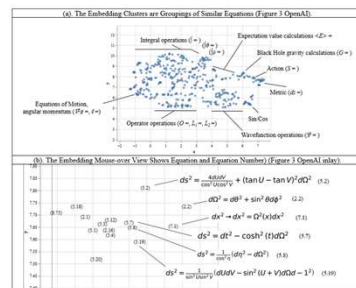
The mathematical embedding visualization of a 476-equation AdS/CFT mathematical ecology [6] is rendered with four standard embedding models (OpenAI, MultiQA, Code Search, MathBERT) in LaTeX and symbolic Python (Figure 2). Interpretatively, there is a distinct structure to the abstracted view of mathematics as embeddings. Most relevant is clustering (not x,y axis values), indicating the grouping of similar kinds of equations, irrespective of order of appearance in the linear progression of a paper. The result is a summary level picture of the kinds mathematics in an ecology, with zoomable equation-image views to examine the mathematics.

Figure 2. AdS/CFT Equation Clusters in Embedding Visualization (LaTeX and SymPy)



Annotated views illustrate (a) how similar groups of equations are grouped in the embedding method and (b) the mouse-over view of equation images by equation number (Figure 3 OpenAI inlay). The initial conversion of equations is to LaTeX, but symbolic Python is a more-readily mobilized computable form of the mathematics for ongoing use e.g. implications for automated equation-solving.

Figure 3. Annotated AdS/CFT Equation Clusters



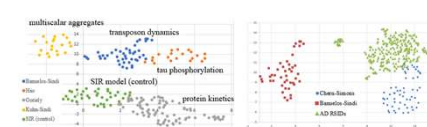
Abbreviations  
AD: Alzheimer's disease; AdS/CFT: anti-de Sitter spacetime conformal field theory; CSF: cerebrospinal fluid; eQTLs: expression quantitative trait loci; GWAS: genome-wide association studies; SNP: single nucleotide polymorphisms; UMAP: uniform manifold approximation and projection (dimensionality reduction)

### Results: Alzheimer's Genomics

Alzheimer's genomics is a whole-human genome-based approach to Alzheimer's disease involving GWAS SNPs, EWAS SNPs, eQTL transcriptomics (expressed RNA), and transposon indels. Whereas ApoE profile was the previous means of assessing Alzheimer's genomic risk, the current understanding includes multiple genomic factors [7]. First, GWAS SNPs may suggest overall risk propensity for the disease. Second, EWAS SNPs indicate which disease genes are actually expressed, confirmed with transcriptome and biomarker assays (Aβ42/40 blood test, CSF). Third, the relation of GWAS-EWAS SNPs on a cis-trans (near-far) correlative basis may be relevant in terms of how different parts of the genome control which genes are expressed. Fourth, transposable elements may be related in activating Alzheimer's disease through genomic insertion-deletion events produced by viruses or other factors.

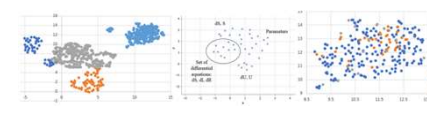
The benefit of the mathematical embedding for descriptive and interventional mathematics in Alzheimer's genomics is that the entirety of a mathscape (set of equations) can be seen in one at-a-glance abstracted and consolidated view. The visualization provides a first-pass view of the mathematics in a paper in the form of equation clusters and mouse-over images. The aggregate view also allows related mathematical ecologies to be compared, and mathematics and data to be investigated in one view as two (ideally corresponding) representations of a system.

Figure 4. Mathematical Ecologies (a) Alzheimer's + SIR Model (control math); (b) Chern-Simons + AD SNPs



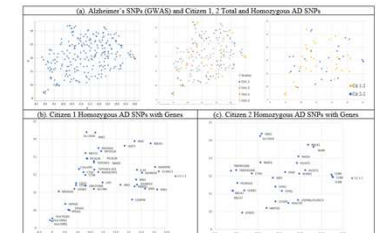
Extending the mathematical embedding from renormalization mathematics to Alzheimer's disease modeling and genomic data, four embedding visualizations for Alzheimer's disease mathscapes are provided, along with a known mathematics control example, the SIR compartmental model [8] (Figure 4a). The lack of overlap is not surprising as the programs target different situations: transposon dynamics, multiscalar aggregates, tau phosphorylation, and protein kinetics [9-12]. The SIR model and the protein kinetics embeddings are closer together as both mathematical ecologies have differential equations as a focal point. The other figure examines mathematics and data together (Figure 4b), suggesting a better model-fit between the Chern-Simons mathematics and the Alzheimer's SNP data than with the transposon dynamics mathematics. Ecosystem analysis is implicated in modeling host-virus interactions in transposon indel activations of Alzheimer's disease. The next steps could include analyzing AdS/CFT mathematical ecologies and the SIR model as applied to multigenic disease indications (Figure 5).

Figure 5. (a) AdS/CFT Mathematical Ecology (b) SIR Mathematics (c) Multi-disease Genomic view: AD, PD, ALS



### Alzheimer's Genomics and SIR Precision Health

Figure 6. Embeddings Visualization of Data: Alzheimer's SNPs applied to Citizen 1, Citizen 2 Precision Health



Alzheimer's disease genomic risk is analyzed for two precision health participants with whole-human genome sequencing (Figure 6). An embedding visualization is performed for all GWAS-linked Alzheimer's disease SNPs and presented for Citizen 1 and Citizen 2's heterozygous (one alternative allele) and homozygous (two alternative alleles) SNPs (a). A gene-level zoom for homozygosity is shown to identify possible pathway-related interventional starting points. Notably, each individual is homozygous for different subsets of genes; Citizen 1 for more immune system related genes (CD33, HLA-DRB1), and Alzheimer's-related clathrin binder (PICALM). Citizen 2 is homozygous for cancer-upregulated membrane proteins (TREM) and cytokine-dependent hematopoietic cell linkers (CLNK). Both are homozygous for the solute carrier protein (SLC24A4) and the intracellular trafficking protein nexin (SNX1). Investigating near-far genomic correlations is enabled by the high-powered mathematical tools of AdS/CFT multiscalar modeling + the UMAP dimensionality-reduction technique in embedding.

### Conclusion

There is an opportunity to deploy AI-based tools to mobilize mathematics as a high-validation data corpus towards broadly humanity-benefiting use cases in global disease-preventing healthy well-being. Future work could elaborate Math Agent uses in an SIR (sustaining, intervening, recovering) model for the societal realization of Precision Health based on two tiers of ongoing circular informational + interventional cycling of the population. Other mathematical discovery use cases of the Math Agent could include synthetic data generation to solve mathematical ecologies, extending equation simplification and math-data model-fit methods, and employing generative AI with episodic memory (per file dating/time-stamping) to assess causal relations in longitudinal personal health dossiers to identify the foundation of pathogenesis. In the short term, genomic variant and eQTL expression data is indicated for practical application to the unresolved challenge of Alzheimer's disease as the top-five human killer with no survivors.

### References

- [1] Nguyen, E., Peil, W., Fiedl, F. et al. (2022). HyenaDNA: Long-Range Genomic Sequence Modeling for Single Nucleotide Resolution. arXiv:2206.15744v1.
- [2] Ballesteros, J. (2022). Large Language Models in Molecular Biology: Deciphering the Language of Biology. Trends Data Science 2 June 2022.
- [3] Houghton, R. (2022). Large Language Models in Molecular Biology: Deciphering the Language of Biology. Trends Data Science 2 June 2022.
- [4] Houghton, R., Houghton, J., & Houghton, J. (2022). Document Clustering for Scientific Literature. arXiv:2206.15744v1.
- [5] Houghton, R. (2022). Document Clustering for Scientific Literature. arXiv:2206.15744v1.
- [6] Wang, J., Houghton, R., & Houghton, J. (2022). Document Clustering for Scientific Literature. arXiv:2206.15744v1.
- [7] Wang, J., Houghton, R., & Houghton, J. (2022). Document Clustering for Scientific Literature. arXiv:2206.15744v1.
- [8] Wang, J., Houghton, R., & Houghton, J. (2022). Document Clustering for Scientific Literature. arXiv:2206.15744v1.
- [9] Wang, J., Houghton, R., & Houghton, J. (2022). Document Clustering for Scientific Literature. arXiv:2206.15744v1.
- [10] Wang, J., Houghton, R., & Houghton, J. (2022). Document Clustering for Scientific Literature. arXiv:2206.15744v1.
- [11] Wang, J., Houghton, R., & Houghton, J. (2022). Document Clustering for Scientific Literature. arXiv:2206.15744v1.
- [12] Wang, J., Houghton, R., & Houghton, J. (2022). Document Clustering for Scientific Literature. arXiv:2206.15744v1.